

# Support vector machines and kernel

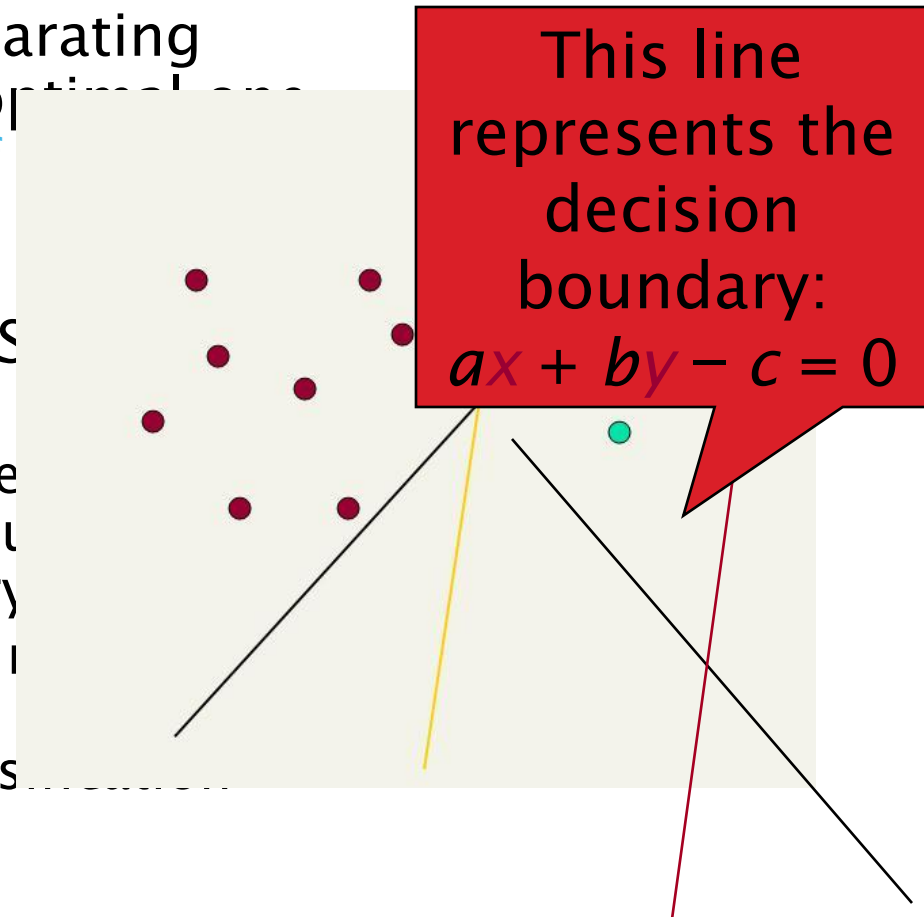


# Text classification: Up until now and today

- ▶ Previously: 3 algorithms for text classification
  - Naive Bayes classifier
  - K Nearest Neighbor classification
    - Simple, expensive at test time, high variance, non-linear
  - Vector space classification using centroids and hyperplanes that split them
    - Simple, linear discriminant classifier; perhaps too simple
      - (or maybe not\*)
- ▶ Today
  - SVMs
  - Some empirical evaluation and comparison
  - Text-specific issues in classification

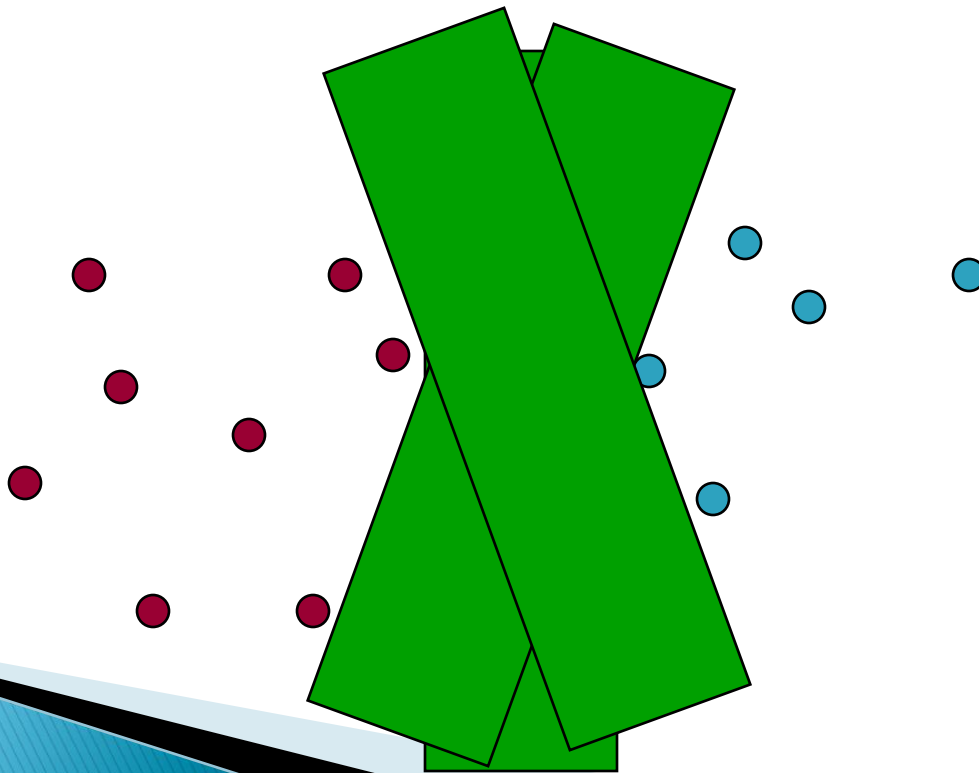
# Linear classifiers: Which Hyperplane?

- ▶ Lots of possible solutions for  $a$ ,  $b$ ,  $c$ .
- ▶ Some methods find a separating hyperplane, but not the optimal one [according to some criterion of goodness]
  - E.g., perceptron
- ▶ Support Vector Machine (SVM) finds an optimal\* solution.
  - Maximizes the distance between the hyperplane and the “difficult” data points (those close to the decision boundary)
  - One intuition: if there are data points very close to the decision surface, the classifier is very uncertain about its class decisions



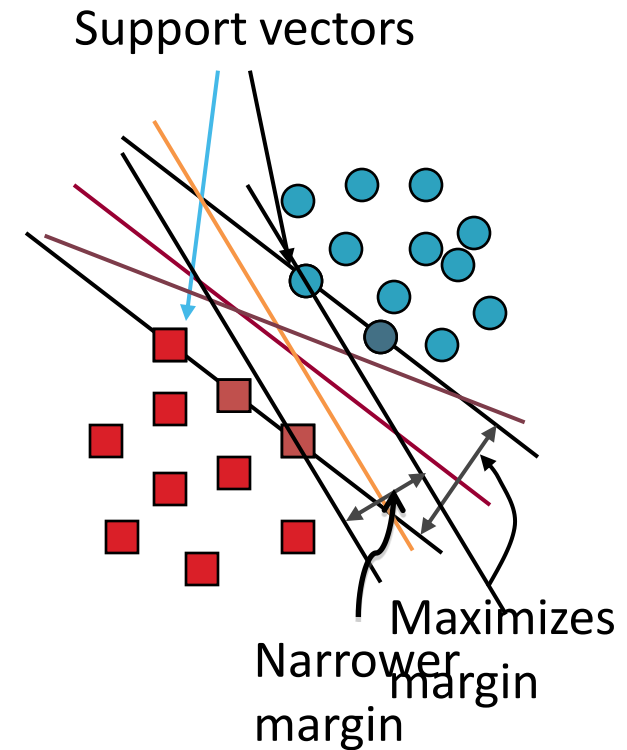
# Another intuition

- ▶ If you have to place a fat separator between classes, you have less choices, and so the capacity of the model has been decreased



# Support Vector Machine (SVM)

- ▶ SVMs maximize the *margin* around the separating hyperplane.
  - A.k.a. large margin classifiers
- ▶ The decision function is fully specified by a subset of training samples, *the support vectors*.
- ▶ Solving SVMs is a *quadratic programming* problem
- ▶ Seen by many as the most successful current text classification method\*



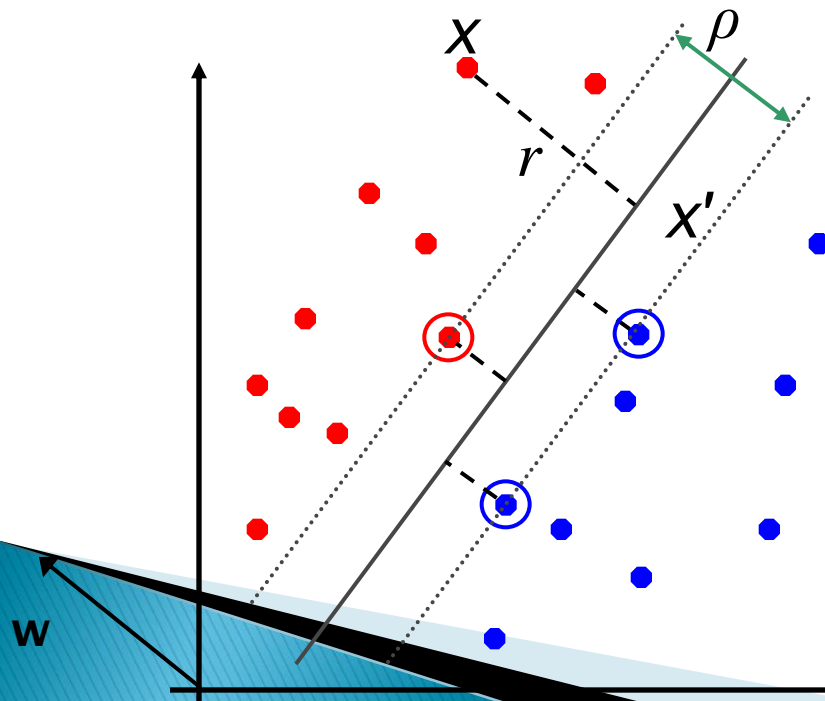
\*but other discriminative methods often perform very similarly

# Maximum Margin: Formalization

- ▶  $\mathbf{w}$ : decision hyperplane normal vector
- ▶  $\mathbf{x}_i$ : data point  $i$
- ▶  $y_i$ : class of data point  $i$  (+1 or -1)      NB: Not 1/0
- ▶ Classifier is:  $f(\mathbf{x}_i) = \text{sign}(\mathbf{w}^T \mathbf{x}_i + b)$
- ▶ Functional margin of  $\mathbf{x}_i$  is:  $y_i (\mathbf{w}^T \mathbf{x}_i + b)$ 
  - But note that we can increase this margin simply by scaling  $\mathbf{w}$ ,  $b$ ....
- ▶ Functional margin of dataset is twice the minimum functional margin for any point
  - The factor of 2 comes from measuring the whole width of the margin

# Geometric Margin

- ▶ Distance from example to the separator is  $y \frac{\mathbf{w}^T \mathbf{x} + b}{\|\mathbf{w}\|}$
- ▶ Examples closest to the hyperplane are **support vectors**.
- ▶ **Margin**  $\rho$  of the separator is the width of separation between support vectors of classes.



Derivation of finding  $r$ :  
Dotted line  $\mathbf{x}' - \mathbf{x}$  is perpendicular to decision boundary so parallel to  $\mathbf{w}$ .

Unit vector is  $\mathbf{w}/\|\mathbf{w}\|$ , so line is  $r\mathbf{w}/\|\mathbf{w}\|$ .

$$\mathbf{x}' = \mathbf{x} - yr\mathbf{w}/\|\mathbf{w}\|.$$

$\mathbf{x}'$  satisfies  $\mathbf{w}^T \mathbf{x}' + b = 0$ .

$$\text{So } \mathbf{w}^T(\mathbf{x} - yr\mathbf{w}/\|\mathbf{w}\|) + b = 0$$

Recall that  $\|\mathbf{w}\| = \sqrt{\mathbf{w}^T \mathbf{w}}$ .

$$\text{So } \mathbf{w}^T \mathbf{x} - yr\|\mathbf{w}\| + b = 0$$

So, solving for  $r$  gives:

$$r = y(\mathbf{w}^T \mathbf{x} + b)/\|\mathbf{w}\|$$

# Linear SVM Mathematically

## The linearly separable case

- ▶ Assume that all data is at least distance 1 from the hyperplane, then the following two constraints follow for a training set  $\{(\mathbf{x}_i, y_i)\}$

$$\mathbf{w}^T \mathbf{x}_i + b \geq 1 \quad \text{if } y_i = 1$$

$$\mathbf{w}^T \mathbf{x}_i + b \leq -1 \quad \text{if } y_i = -1$$

- ▶ For support vectors, the inequality becomes an equality
- ▶ Then, since each example's distance from the hyperplane is

$$r = y \frac{\mathbf{w}^T \mathbf{x} + b}{\|\mathbf{w}\|}$$

- ▶ The margin is:

$$r = \frac{2}{\|\mathbf{w}\|}$$



# Linear Support Vector Machine (SVM)

▶ **Hyperplane**

$$w^T x + b = 0$$

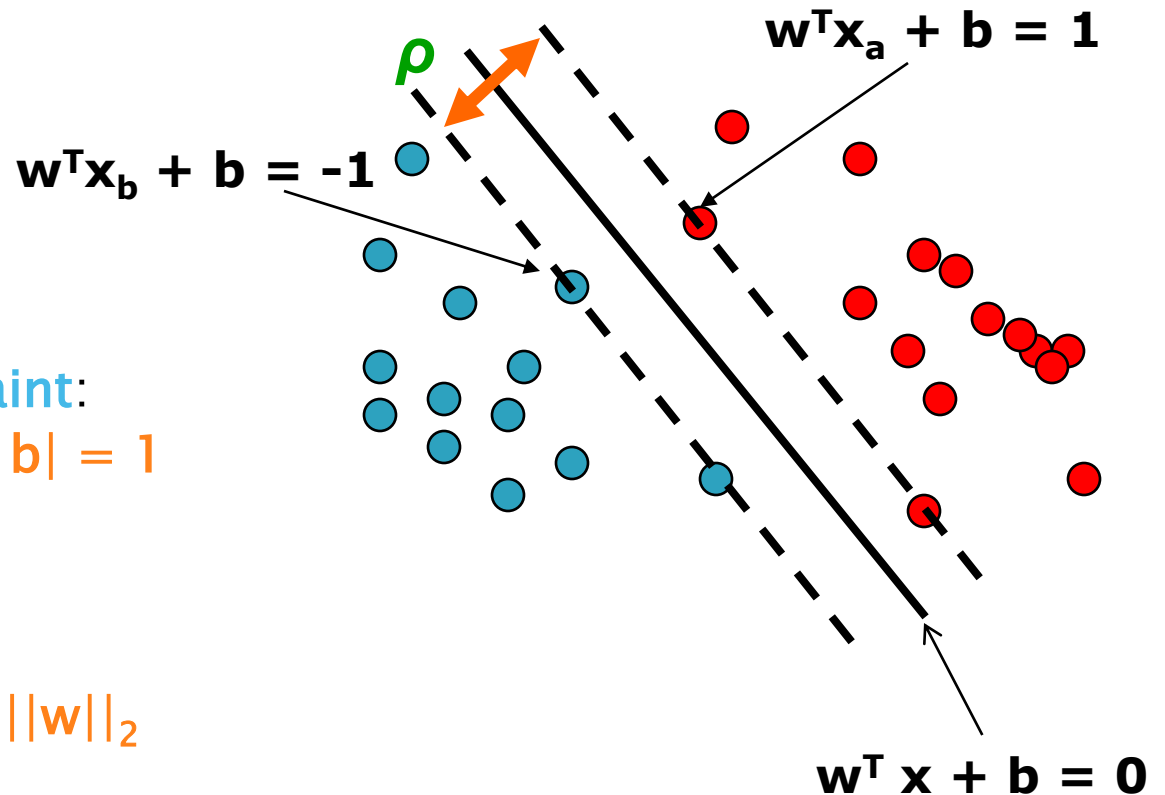
▶ **Extra scale constraint:**

$$\min_{i=1, \dots, n} |w^T x_i + b| = 1$$

▶ **This implies:**

$$w^T(x_a - x_b) = 2$$

$$\rho = \|x_a - x_b\|_2 = 2 / \|w\|_2$$



# Linear SVMs Mathematically (cont.)

- ▶ Then we can formulate the *quadratic optimization problem*:

Find  $\mathbf{w}$  and  $b$  such that

$$r = \frac{2}{\|\mathbf{w}\|} \quad \text{is maximized; and for all } \{(\mathbf{x}_i, y_i)\}$$
$$\mathbf{w}^T \mathbf{x}_i + b \geq 1 \text{ if } y_i=1; \quad \mathbf{w}^T \mathbf{x}_i + b \leq -1 \text{ if } y_i = -1$$

- ▶ A better formulation ( $\min \|\mathbf{w}\| = \max 1 / \|\mathbf{w}\|$ ):

Find  $\mathbf{w}$  and  $b$  such that

$\Phi(\mathbf{w}) = 1/2 \mathbf{w}^T \mathbf{w}$  is minimized;

and for all  $\{(\mathbf{x}_i, y_i)\}$ :  $y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1$

# Linear SVMs: Summary

- ▶ The classifier is a *separating hyperplane*.
- ▶ The most “important” training points are the support vectors; they define the hyperplane.
- ▶ Quadratic optimization algorithms can identify which training points  $\mathbf{x}_i$  are support vectors with non-zero Lagrangian multipliers  $\alpha_i$ .
- ▶ Both in the dual formulation of the problem and in the solution, training points appear only inside inner products:

Find  $\alpha_1 \dots \alpha_N$  such that

$\mathbf{Q}(\boldsymbol{\alpha}) = \sum \alpha_i - \frac{1}{2} \sum \sum \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j$  is maximized and

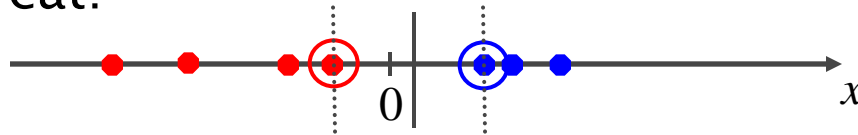
(1)  $\sum \alpha_i y_i = 0$

(2)  $0 \leq \alpha_i \leq C$  for all  $\alpha_i$

$$f(\mathbf{x}) = \sum \alpha_i y_i \mathbf{x}_i^T \mathbf{x} + b$$

# Non-linear SVMs

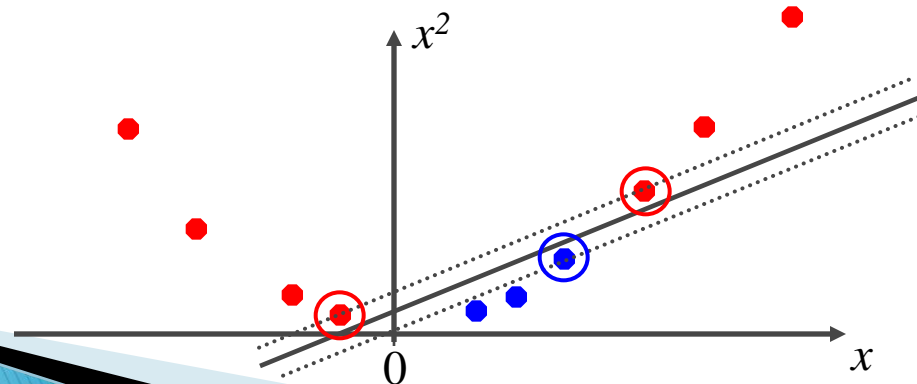
- ▶ Datasets that are linearly separable (with some noise) work out great:



- ▶ But what are we going to do if the dataset is just too hard?

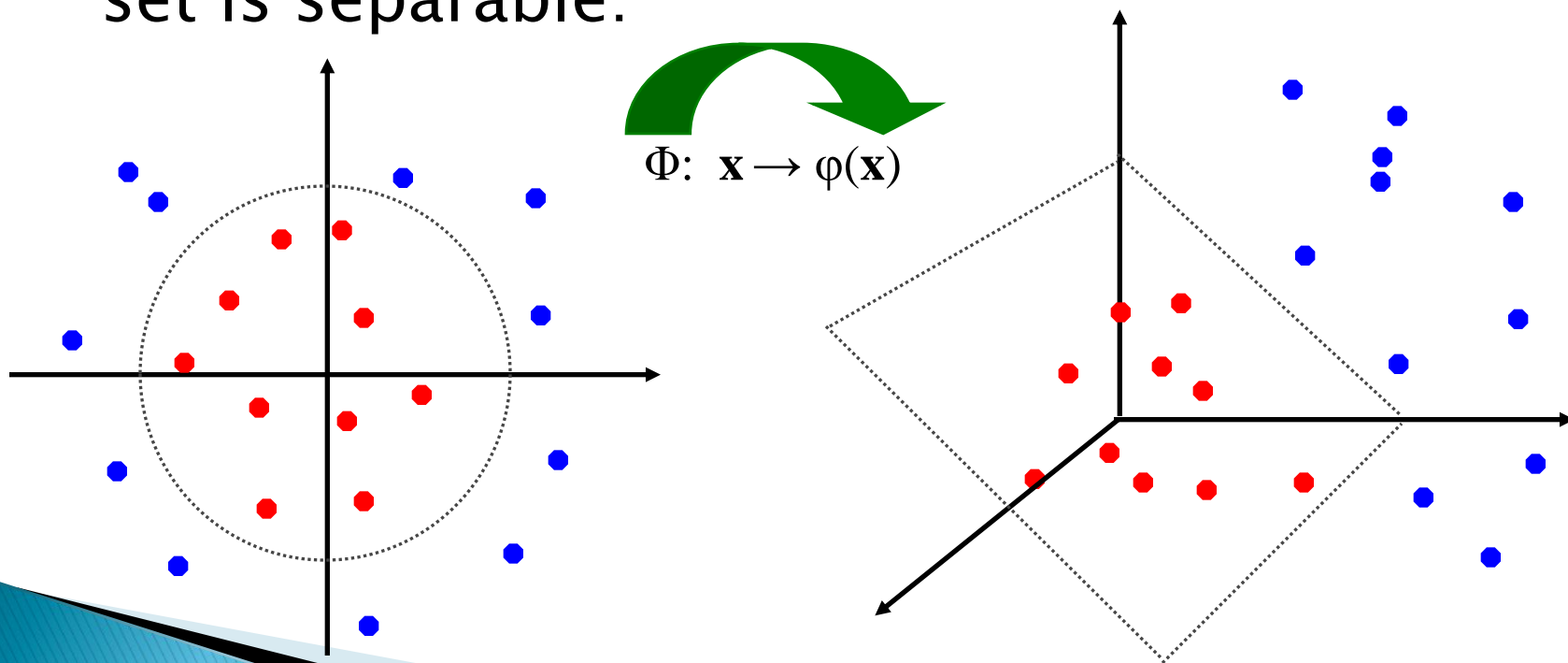


- ▶ How about ... mapping data to a higher-dimensional space:



# Non-linear SVMs: Feature spaces

- ▶ General idea: the original feature space can always be mapped to some higher-dimensional feature space where the training set is separable:



# The “Kernel Trick”

- ▶ The linear classifier relies on an inner product between vectors  
 $K(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i^T \mathbf{x}_j$
- ▶ If every datapoint is mapped into high-dimensional space via some transformation  $\Phi: \mathbf{x} \rightarrow \phi(\mathbf{x})$ , the inner product becomes:

$$K(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$$

- ▶ A *kernel function* is some function that corresponds to an inner product in some expanded feature space.
- ▶ Example:

2-dimensional vectors  $\mathbf{x} = [x_1 \ x_2]$ ; let  $K(\mathbf{x}_i, \mathbf{x}_j) = (1 + \mathbf{x}_i^T \mathbf{x}_j)^2$ ,

Need to show that  $K(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$ :

$$K(\mathbf{x}_i, \mathbf{x}_j) = (1 + \mathbf{x}_i^T \mathbf{x}_j)^2 = 1 + x_{i1}^2 x_{j1}^2 + 2 x_{i1} x_{j1} x_{i2} x_{j2} + x_{i2}^2 x_{j2}^2 + 2 x_{i1} x_{j1} + 2 x_{i2} x_{j2} =$$

$$= [1 \ x_{i1}^2 \ \sqrt{2} x_{i1} x_{i2} \ x_{i2}^2 \ \sqrt{2} x_{i1} \ \sqrt{2} x_{i2}]^T [1 \ x_{j1}^2 \ \sqrt{2} x_{j1} x_{j2} \ x_{j2}^2 \ \sqrt{2} x_{j1} \ \sqrt{2} x_{j2}]$$

$$= \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j) \quad \text{where } \phi(\mathbf{x}) = [1 \ x_1^2 \ \sqrt{2} x_1 x_2 \ x_2^2 \ \sqrt{2} x_1 \ \sqrt{2} x_2]$$

# Kernels

- ▶ Why use kernels?
  - Make non-separable problem separable.
  - Map data into better representational space
- ▶ Common kernels
  - Linear
  - Polynomial  $K(\mathbf{x}, \mathbf{z}) = (1 + \mathbf{x}^T \mathbf{z})^d$ 
    - Gives feature conjunctions
  - Radial basis function (infinite dimensional space)

▶ Hav clas  $K(\mathbf{x}_i, \mathbf{x}_j) = e^{-\|\mathbf{x}_i - \mathbf{x}_j\|^2 / 2\sigma^2}$

# Per class evaluation measures

- ▶ Recall: Fraction of docs in class  $i$  classified correctly:  $\frac{c_{ii}}{\sum_j c_{ij}}$
- ▶ Precision: Fraction of docs assigned class  $i$  that are actually about class  $i$ :  $\frac{c_{ii}}{\sum_j c_{ji}}$
- ▶ Accuracy: (1 – error rate) Fraction of docs classified correctly:  $\frac{\sum_i c_{ii}}{\sum_j \sum_i c_{ij}}$



# Micro- vs. Macro-Averaging

- ▶ If we have more than one class, how do we combine multiple performance measures into one quantity?
- ▶ Macroaveraging: Compute performance for each class, then average.
- ▶ Microaveraging: Collect decisions for all classes, compute contingency table, evaluate.

# Micro- vs. Macro-Averaging: Example

Class 1

	Truth: yes	Truth: no
Classifier: yes	10	10
Classifier: no	10	970

Class 2

	Truth: yes	Truth: no
Classifier: yes	90	10
Classifier: no	10	890

Micro Ave. Table

	Truth: yes	Truth: no
Classifier: yes	100	20
Classifier: no	20	1860

- Macroaveraged precision:  $(0.5 + 0.9)/2 = 0.7$
- Microaveraged precision:  $100/120 = .83$
- Microaveraged score is dominated by score on common classes

# The Real World

P. Jackson and I. Moulinier. 2002. *Natural Language Processing for Online Applications*

- ▶ “There is no question concerning the commercial value of being able to classify documents automatically by content. There are myriad potential applications of such a capability for corporate intranets, government departments, and Internet publishers”
- ▶ “Understanding the data is one of the keys to successful categorization, yet this is an area in which most categorization tool vendors are extremely weak. Many of the ‘one size fits all’ tools on the market have not been tested on a wide range of content types.”

# The Real World

- ▶ Gee, I'm building a text classifier for real, now!
- ▶ What should I do?
  
- ▶ How much training data do you have?
  - None
  - Very little
  - Quite a lot
  - A huge amount and its growing

# Manually written rules

- ▶ No training data, adequate editorial staff?
- ▶ Never forget the hand-written rules solution!
  - If (wheat or grain) and not (whole or bread) then
    - Categorize as grain
- ▶ In practice, rules get a lot bigger than this
  - Can also be phrased using tf or tf.idf weights
- ▶ With careful crafting (human tuning on development data) performance is high:
  - Construe: 94% recall, 84% precision over 675 categories (Hayes and Weinstein 1990)
- ▶ Amount of work required is huge
  - Estimate 2 days per class ... plus maintenance

# Very little data?

- ▶ If you're just doing supervised classification, you should stick to something high bias
  - There are theoretical results that Naïve Bayes should do well in such circumstances (Ng and Jordan 2002 NIPS)
- ▶ The interesting theoretical answer is to explore semi-supervised training methods:
  - Bootstrapping, EM over unlabeled documents, ...
- ▶ The practical answer is to get more labeled data as soon as you can
  - How can you insert yourself into a process where humans will be willing to label data for you??

# A reasonable amount of data?

- ▶ Perfect!
- ▶ We can use all our clever classifiers
- ▶ Roll out the SVM!
  
- ▶ But if you are using an SVM/NB etc., you should probably be prepared with the “hybrid” solution where there is a Boolean overlay
  - Or else to use user–interpretable Boolean–like models like decision trees
  - Users like to hack, and management likes to be able to implement quick fixes immediately

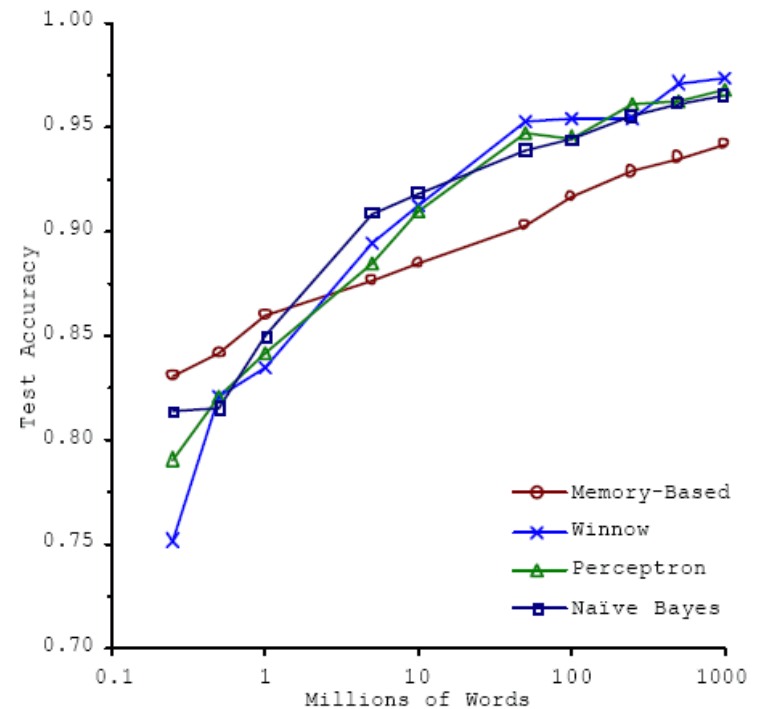
# A huge amount of data?

- ▶ This is great in theory for doing accurate classification...
- ▶ But it could easily mean that expensive methods like SVMs (train time) or kNN (test time) are quite impractical
- ▶ Naïve Bayes can come back into its own again!
  - Or other advanced methods with linear training/test complexity like regularized logistic regression (though much more expensive to train)



# Accuracy as a function of data size

- ▶ With enough data the choice of classifier may not matter much, and the best choice may be unclear
  - Data: Brill and Banko on context-sensitive spelling correction
- ▶ But the fact that you have to keep doubling your data to improve performance is a little unpleasant



# How many categories?

- ▶ A few (well separated ones)?
  - Easy!
- ▶ A zillion closely related ones?
  - Think: Yahoo! Directory, Library of Congress classification, legal applications
  - Quickly gets difficult!
    - Classifier combination is always a useful technique
      - Voting, bagging, or boosting multiple classifiers
    - Much literature on hierarchical classification
      - Mileage fairly unclear, but helps a bit (Tie-Yan Liu et al. 2005)
    - May need a hybrid automatic/manual solution

# How can one tweak performance?

- ▶ Aim to exploit any domain-specific useful features that give special meanings or that zone the data
  - E.g., an author byline or mail headers
- ▶ Aim to collapse things that would be treated as different but shouldn't be.
  - E.g., part numbers, chemical formulas
- ▶ Does putting in “hacks” help?
  - You bet!
    - Feature design and non-linear weighting is *very* important in the performance of real-world systems

# Upweighting

- ▶ You can get a lot of value by differentially weighting contributions from different document zones:
- ▶ That is, you count as two instances of a word when you see it in, say, the abstract
  - Upweighting title words helps (Cohen & Singer 1996)
    - Doubling the weighting on the title words is a good rule of thumb
  - Upweighting the first sentence of each paragraph helps (Murata, 1999)
  - Upweighting sentences that contain title words helps (Ko *et al*, 2002)

# Two techniques for zones

1. Have a completely separate set of features/parameters for different zones like the title
  2. Use the same features (pooling/tying their parameters) across zones, but upweight the contribution of different zones
- ▶ Commonly the second method is more successful: it costs you nothing in terms of sparsifying the data, but can give a very useful performance boost
    - Which is best is a contingent fact about the data

# Text Summarization techniques in text classification

- ▶ Text Summarization: Process of extracting key pieces from text, normally by features on sentences reflecting position and content
- ▶ Much of this work can be used to suggest weightings for terms in text categorization
  - See: Kolcz, Prabakarmurthi, and Kalita, CIKM 2001: Summarization as feature selection for text categorization
  - Categorizing purely with title,
  - Categorizing with first paragraph only
  - Categorizing with paragraph with most keywords
  - Categorizing with first and last paragraphs, etc.

# Does stemming/lowercasing/... help?

- ▶ As always, it's hard to tell, and empirical evaluation is normally the gold standard
- ▶ But note that the role of tools like stemming is rather different for TextCat vs. IR:
  - For IR, you often want to collapse forms of the verb *oxygenate* and *oxygenation*, since all of those documents will be relevant to a query for *oxygenation*
  - For TextCat, with sufficient training data, stemming *does no good*. It only helps in compensating for data sparseness (which can be severe in TextCat applications). *Overly aggressive stemming can easily degrade performance.*

# Measuring Classification Figures of Merit

- ▶ Not just accuracy; in the real world, there are economic measures:
  - Your choices are:
    - Do no classification
      - That has a cost (hard to compute)
    - Do it all manually
      - Has an easy-to-compute cost if doing it like that now
    - Do it all with an automatic classifier
      - Mistakes have a cost
    - Do it with a combination of automatic classification and manual review of uncertain/difficult/"new" cases
  - Commonly the last method is most cost efficient and is adopted



# A common problem: Concept Drift

- ▶ Categories change over time
- ▶ Example: “president of the united states”
  - 1999: clinton is great feature
  - 2010: clinton is bad feature
- ▶ One measure of a text classification system is how well it protects against concept drift.
  - Favors simpler models like Naïve Bayes
- ▶ Feature selection: can be bad in protecting against concept drift

# Summary

- ▶ Support vector machines (SVM)
    - Choose hyperplane based on support vectors
      - Support vector = “critical” point close to decision boundary
    - (Degree-1) SVMs are linear classifiers.
    - Kernels: powerful and elegant way to define similarity metric
    - Perhaps best performing text classifier
      - But there are other methods that perform about as well as SVM, such as regularized logistic regression (Zhang & Oles 2001)
    - Partly popular due to availability of good software
      - SVMlight is accurate and fast – and free (for research)
      - Now lots of good software: libsvm, TinySVM, ....
  - ▶ Comparative evaluation of methods
- Real world: exploit domain specific structure!